

Reproducible Experiment Platform

Alexander Baranov, Egor Khairullin, Tatiana Likhomanenko, Alexey Rogozhnikov, Andrey Ustyuzhanin

The 21st International Conference on Computing in High Energy and Nuclear Physics, Okinawa, 2015



github.com/yandex/REP



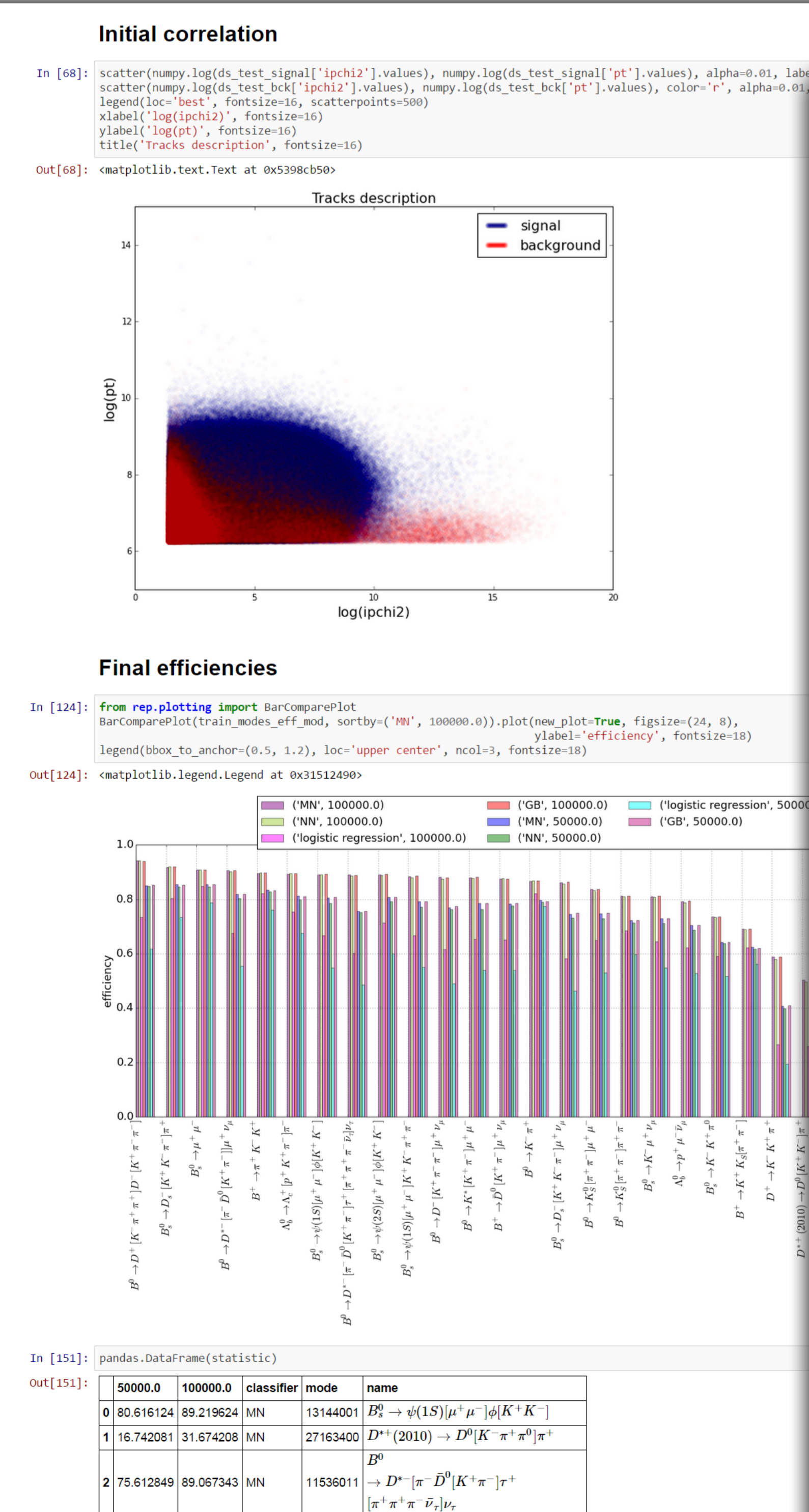
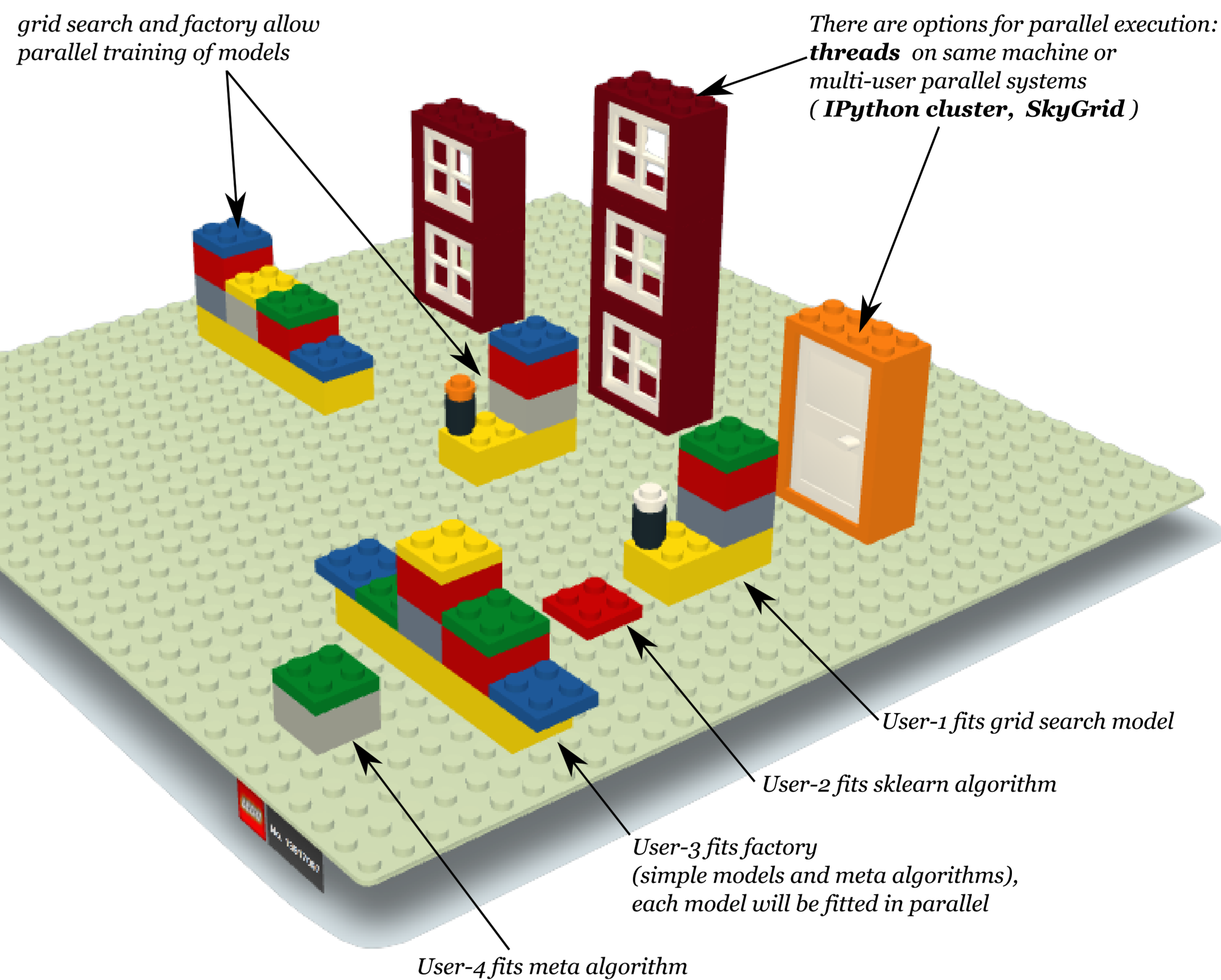
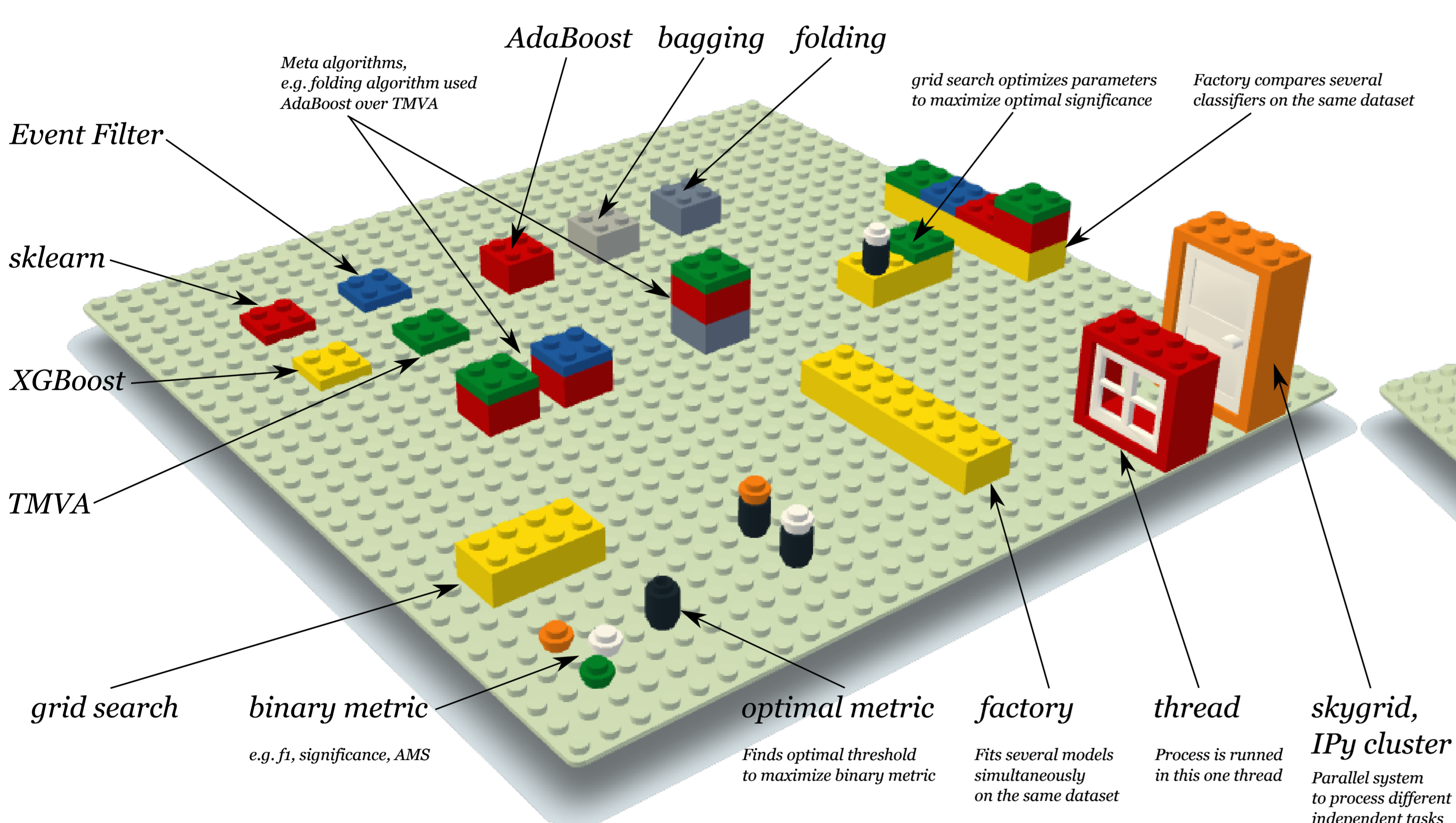
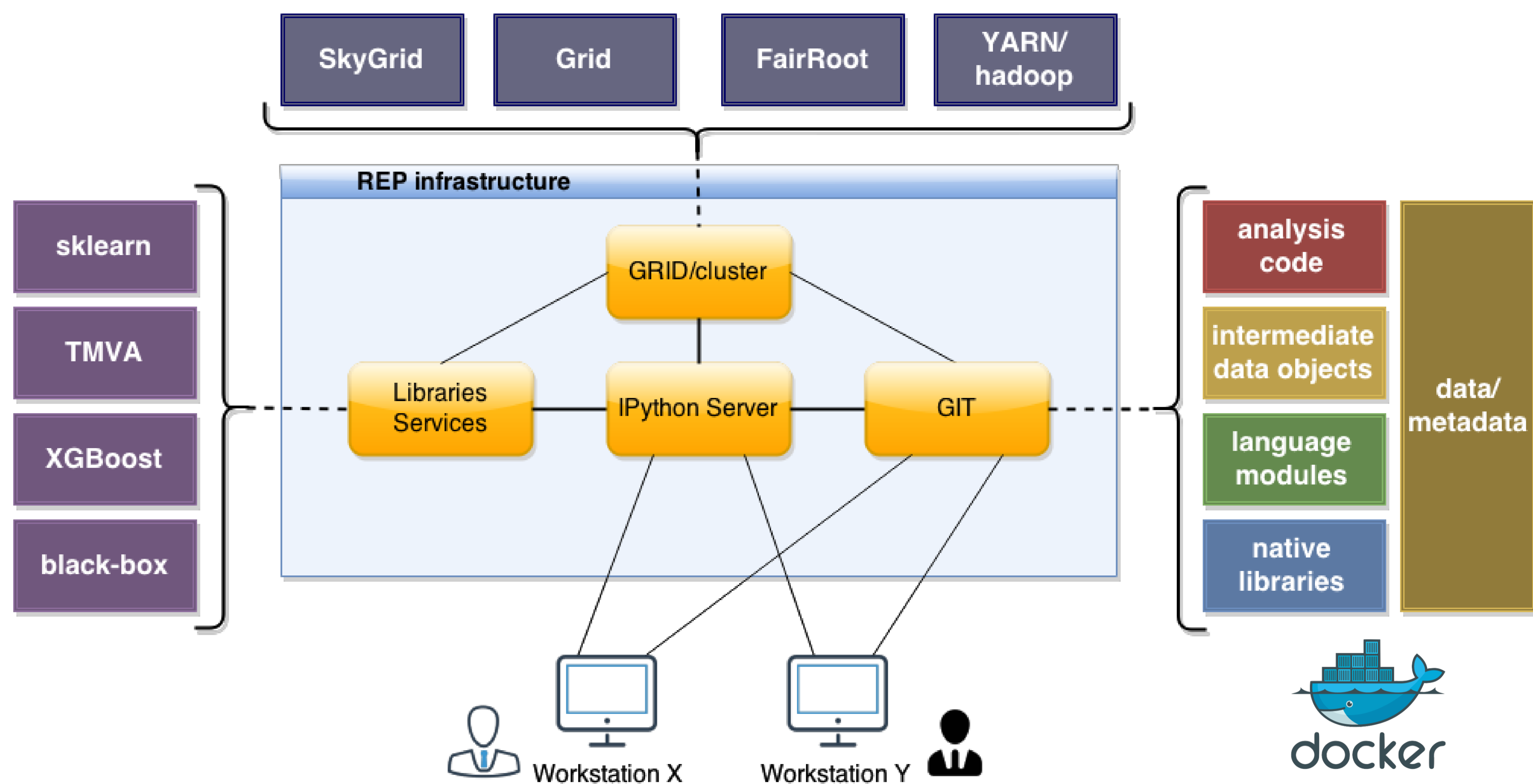
“Non-reproducible single occurrences are of no significance to science.” – Karl Popper

Plan for paper:

- Describe theory
- Some **magic** happens
- Conclusions



Great! Finally I found an appropriate journal to publish my discovery!



Viva la Factory

Variables needed for analysis

```

In [1]: loaded_variables = ["FlightDistance", "FlightDistanceError", "IP", "VertexChi2", "pt", "p0", "mass"]
train_variables = ["FL: FlightDistance/FlightDistanceError", "IP", "VertexChi2", "pt", "p0"]
plot_variables = train_variables + ["mass"]

```

Loading data

```

In [2]: import numpy, pandas
from rep.utils import train_test_split

In [3]: sig_data = pandas.read_csv('toy_datasets/toyMC_sig_mass.csv', sep='t', usecols=loaded_variables)
bck_data = pandas.read_csv('toy_datasets/toyMC_bck_mass.csv', sep='t', usecols=loaded_variables)

labels = numpy.array([1] * len(sig_data) + [0] * len(bck_data))
data = pandas.concat([sig_data, bck_data])

# Get train and test data
train_data, test_data, train_labels, test_labels = train_test_split(data, labels, train_size=0.7)

```

Factory of different models

```

In [4]: from rep.metalm import ClassifiersFactory
from rep.estimators import EventFilterClassifier, TMVAClassifier, sklearnClassifier, XGBoostClassifier
from sklearn.ensemble import AdaBoostClassifier

In [5]: fac = factory(name, cl in estimators.items()):
# there are different ways to add classifiers to Factory:
factory['ef'] = EventFilterClassifier(connection='test_connection', features=train_variables, iterations=200)
factory.add_classifier('tmva', TMVAClassifier(NTrees=50, features=train_variables[:5], Shrinkage=0.1))
factory.add_classifier('ada', AdaBoostClassifier(n_estimators=10))
factory['xgb'] = XGBoostClassifier(features=train_variables[2:6])
factory.fit(train_data, train_labels)

```

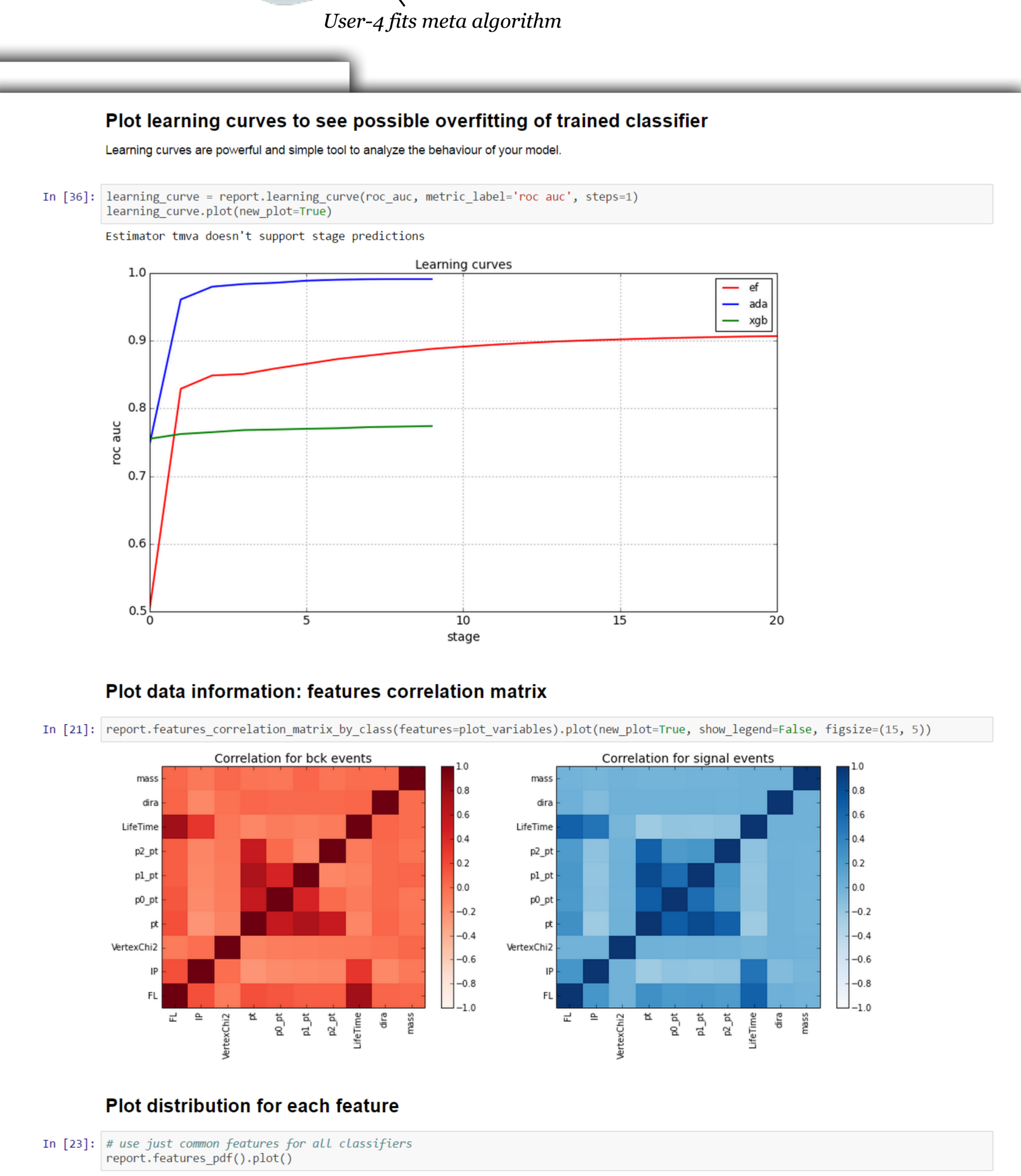
Everybody loves plots!

```

In [7]: report = factory.test_on(test_data, test_labels)

In [9]: report.roc().plot(xlim=(0.5, 1))

```



REP was used for:

- Data popularity estimator github.com/hushchyn-mikhail/DataPopularity
- Uboost-like algorithms github.com/anaderi/lhcb_trigger_ml
- Bs -> mu mu analysis
- LHCb topological trigger optimization



References:
 • <http://buildwithchrome.com>
 • <http://jir.com>